

# Simple and Efficient Parallelization for Probabilistic Temporal Tensor Factorization

Guangxi Li<sup>1</sup>, Zenglin Xu<sup>1</sup>, Linnan Wang<sup>1</sup>, Jinmian Ye<sup>1</sup>, Irwin King<sup>2</sup>, Michael Lyu<sup>2</sup>

<sup>1</sup> Big Data Research Center, School of Computer Science and Engineering, University of Electronic Science and Technology of China, Chengdu, Sichuan, China;

gxli@std.uestc.edu.cn; {zenglin, wangnan318}@gmail.com; me@ay27.com

<sup>2</sup> Department of Computer Science and Technology, The Chinese University of Hong Kong Shatian, N.T., Hong Kong; {king, lyu}@cse.cuhk.edu.hk

**Abstract**—*Probabilistic Temporal Tensor Factorization (PTTF)* is an effective algorithm to model the temporal tensor data. It leverages a time constraint to capture the evolving properties of tensor data. Nowadays the exploding dataset demands a large scale PTTF analysis, and a parallel solution is critical to accommodate the trend. Whereas, the parallelization of PTTF still remains unexplored. In this paper, we propose a simple yet efficient *Parallel Probabilistic Temporal Tensor Factorization*, referred to as  $P^2T^2F$ , to provide a scalable PTTF solution.  $P^2T^2F$  is fundamentally disparate from existing parallel tensor factorizations by considering the probabilistic decomposition and the temporal effects of tensor data. It adopts a new tensor data split strategy to subdivide a large tensor into independent sub-tensors, the computation of which is inherently parallel. We train  $P^2T^2F$  with an efficient algorithm of stochastic Alternating Direction Method of Multipliers, and show that the convergence is guaranteed. Experiments on several real-world tensor datasets demonstrate that  $P^2T^2F$  is a highly effective and efficiently scalable algorithm dedicated for large scale probabilistic temporal tensor analysis.

## I. INTRODUCTION

Recent developments of tensor decomposition have great impacts on signal processing [1], computer vision [2], numerical analysis [3], social network analysis [4], recommendation systems [5] and etc. A comprehensive overview can be found from the survey paper by [6]. In particular, automatic recommendation systems significantly benefit from tensor decomposition as it effectively extracts hidden patterns from the multi-way data.

Various tensor decomposition methods have been proposed. The CANDECOMP/PARAFAC decomposition, shorted as CP decomposition, is a direct extension of low-rank matrix decomposition to tensors; and it can be regarded as a special case of Tucker Decomposition by adding a super-diagonal constraint on the core tensor [7]. This method, however, fails to consider the fact that the real relational data is evolving over time and exhibits strong temporal patterns, especially in recommendation systems. To resolve this issue, *Probabilistic Temporal Tensor Factorization* (PTTF) [8], inspired by probabilistic latent factor models [9], [10], has been proposed by incorporating a time constraint. In contrast with Multi-HDP [11], Probabilistic Non-negative Tensor Factorization [12] and

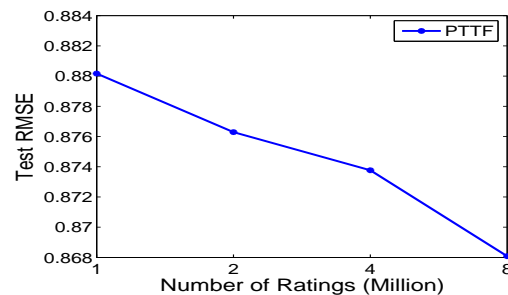


Fig. 1. The RMSE curve for PTTF on MovieLens. The prediction error reduces as the number of ratings grows.

Probabilistic Polyadic Factorization [13], PTTF is the only one capturing the temporal effects of tensor data.

The era of big data has also witnessed the explosion of tensor datasets, while the large scale PTTF analysis is important to accommodate the increasing datasets. Figure 1 demonstrates PTTF achieves better performance as the tensor size increases on the MovieLens data (Table I). The result directly sheds light on the necessity of a parallel PTTF solution. Nevertheless, there is a huge gap to be filled.

In this paper, we present *Parallel Probabilistic Temporal Tensor Factorization* ( $P^2T^2F$ ) dedicated for large-scale temporal tensor factorization problems. The core concept of  $P^2T^2F$  is to reduce each sequential operation on a large tensor into a set of independent operations on smaller sub-tensors for parallel executions, while still retains the ability to model temporal effects of PTTF. In general, the main contributions of  $P^2T^2F$  are as follows:

- $P^2T^2F$  allows parallel solutions to probabilistic tensor decomposition with temporal effects and thus makes PTTF model scalable. We demonstrate a new parallelization scheme in  $P^2T^2F$  to divide the large-scale problem into several sub-problems for concurrent executions.
- In  $P^2T^2F$ , we also design a novel stochastic learning algorithm for parallel ADMM framework to improve the PTTF model. Specifically, this algorithm calculates the latent feature factors using a substitutive objective function that is convex and can be viewed as an upper

bound of the original problem.

- The convergence of P<sup>2</sup>T<sup>2</sup>F is theoretically guaranteed.

## II. RELATED WORK

Matrix or tensor factorization methods are useful tools in recommendation systems. One prominent representative factor-based method for recommendation systems is *Probabilistic Matrix Factorization* (PMF) [9], the latent factors of which can be learned by maximum likelihood estimation. Temporal modeling has been greatly ignored in the community of collaborative filtering until the timeSVD++ algorithm is proposed in [14]. This method demonstrates that the latent features consist of components evolving over time; and such features effectively capture local changes of user preferences. To extend the method to tensors, [8] proposes the PTTF model to capture the global effect of time shared among users and items.

The increasing demand of modeling data scalability incubates several parallel models for PMF problems, such as PPMF [15], Hogwild [16] and DSGD [17]. However, compared to large-scale matrix factorization [18], there are fewer works devoted to large-scale tensors. In general, the existing large-scale tensor methods can be categorized into two classes. The first one consists in exploiting sparseness of tensors. For example, the GigaTensor algorithm in [19] and DFacTo method in [20] intend to minimize the number of floating point operations and to handle the size of intermediate data to avoid the intermediate data explosion problem, respectively. The other class of methods consists in distributing the computation load to a number of workers [21], [22], [23]. Unfortunately, these algorithms do not concentrate on modeling the temporal effects of tensor data.

ADMM is an effective framework for accelerating the optimization of tensor factorization [24], [25], [23]. However, they do not target for improving the scalability of tensor factorization algorithms; it also neglects to model the temporal effects in dynamic systems. For example, the parallel ADMM algorithm proposed in [23] computes tensor decomposition mode by mode; and it is hard to be deployed in online training or distributed training. As a response, we propose P<sup>2</sup>T<sup>2</sup>F that adapts for the online training or the distributed training. It also works in various computing environments, such as multi-core computers or clusters.

## III. PRELIMINARIES

### A. Notations

A tensor is a multi-dimensional array that generalizes a vectors (1-dimensional tensor) and a matrix (2-dimensional tensor) to higher order. Like rows and columns in a matrix, an N-dimensional tensor has N modes whose lengths are  $I_1$  through  $I_N$ , respectively. By convention, vectors and matrices are denoted by boldface lowercase letters or uppercase letters with a subscript, e.g.,  $\mathbf{a}$  or  $\mathbf{A}_i$ , and boldface capital letters, e.g.,  $\mathbf{A}$ , respectively. We denote higher-order tensors (order three or higher) by boldface Euler script letters, e.g.,  $\mathcal{X}$ . We also denote the entry of a tensor by the symbolic name of

tensor with its indices in subscript. For example, the  $(i_1, i_2)$ th entry of  $\mathbf{A}$  is denoted by  $a_{i_1 i_2}$ , and the  $(i_1, \dots, i_N)$ th entry of  $\mathcal{X}$  is denoted by  $x_{i_1, \dots, i_N}$ .

### B. Tensor Decomposition

There are several ways to define tensor decomposition [6], [26], [27]. Our definition is based on CP (CANDECOMP/PARAFAC) decomposition, which is one of the most popular decomposition methods. Details about CP decomposition can be found in [6]. For ease of presentation, we only derive our model in the third-order case, but it can be easily generalized to N-way tensor.

Let  $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$  be a third order tensor with observable entries  $\{x_{ijk} | (i, j, k) \in \Omega\}$ , we hope to find the factor matrices  $\{\mathbf{A} \in \mathbb{R}^{I \times R}, \mathbf{B} \in \mathbb{R}^{J \times R}, \mathbf{C} \in \mathbb{R}^{K \times R}\}$  by minimizing the following loss function:

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \frac{1}{2} \sum_{(i,j,k) \in \Omega} (x_{ijk} - \langle \mathbf{A}_i, \mathbf{B}_j, \mathbf{C}_k \rangle)^2 + \frac{\lambda}{2} (\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2 + \|\mathbf{C}\|_F^2), \quad (1)$$

where  $\langle \mathbf{A}_i, \mathbf{B}_j, \mathbf{C}_k \rangle \equiv \sum_{r=1}^R a_{ir} b_{jr} c_{kr}$  denotes the inner product of three  $R$ -dimensional vectors,  $\mathbf{A}_i$  denotes the  $i$ th row of  $\mathbf{A}$ , so does  $\mathbf{B}_j$  and  $\mathbf{C}_k$ . If we solve this model with *Stochastic Gradient Descent* (SGD), a drastic simplification [28], we can get a SGD-based CP decomposition model. The stochastic process depends on the examples randomly drawn at each iteration.

### C. Probabilistic Temporal Tensor Factorization

*Probabilistic Temporal Tensor Factorization* (PTTF) model can be considered as the extension of PMF model [9] by adding a specially-constrained time dimension. For the third order tensor  $\mathcal{X}$  in (1), if the third dimension denotes the time corresponding to the factor matrix  $\mathbf{C}$ , we assume the following conditional prior for  $\mathbf{C}$  [8]:

$$\mathbf{C}_k \sim \mathcal{N}(\mathbf{C}_{k-1}, \sigma_C^2 \mathbf{I}_R), \quad k = 1, \dots, K.$$

For the initial time feature vector  $\mathbf{C}_0$ , we assume

$$\mathbf{C}_0 \sim \mathcal{N}(\boldsymbol{\mu}_C, \sigma_0^2 \mathbf{I}_R),$$

where  $\boldsymbol{\mu}_C$  denotes a 1-by- $R$  row vector and  $\mathbf{I}_R$  denotes a  $R$ -by- $R$  identity matrix. The PTTF model can be expressed to minimizing the following regularized sum of squared errors (the proof can be seen in [8]):

$$\frac{1}{2} \sum_{(i,j,k) \in \Omega} (x_{ijk} - \langle \mathbf{A}_i, \mathbf{B}_j, \mathbf{C}_k \rangle)^2 + \frac{\lambda_0}{2} \|\mathbf{C}_0 - \boldsymbol{\mu}_C\|_2^2 + \frac{\lambda_A}{2} \|\mathbf{A}\|_F^2 + \frac{\lambda_B}{2} \|\mathbf{B}\|_F^2 + \frac{\lambda_C}{2} \sum_{k=1}^K \|\mathbf{C}_k - \mathbf{C}_{k-1}\|_2^2, \quad (2)$$

where  $\lambda_A, \lambda_B, \lambda_C, \lambda_0$  denote regularization parameters.

Obviously, we can also get a SGD-based PTTF model using SGD method, but it can not be easily parallelized because of the special constraint on the time dimension. In

---

**Algorithm 1** Tensor Data Split

---

**Input:** Tensor  $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ ,  $P$ .

**Output:**  $\mathcal{X}^p \in \mathbb{R}^{\frac{I}{P} \times J \times K}$ .

- 1) initialize  $I_{first} = 0$ ;
  - 2) **for**  $p = 1, \dots, P$  **do**
  - 3)  $I_{last} = \lfloor \frac{I}{P} \rfloor$ ;
  - 4)  $\mathcal{X}^p \leftarrow \mathcal{X}_{(I_{first}+1:I_{last}), :, :}$ ;
  - 5)  $I_{first} = I_{last}$ ;
  - 6) **return**  $\mathcal{X}^p$ .
- 

contrast, ADMM framework has the properties of flexibility and tractability and, the most important, can be naturally used to design parallel or distributed learning algorithms for large-scale problems. So, in the next section, we will address the minimization problem in (2) in the framework of ADMM. In recent years, ADMM has occupied more and more attention and wide range of applications such as matrix completion [15], [29] and compressive sensing [30], but as far as we know, few works have been proposed to use parallel ADMM framework for probabilistic temporal tensor decomposition problems.

#### IV. PARALLEL PROBABILISTIC TEMPORAL TENSOR FACTORIZATION

In this section, we elaborate *Parallel Probabilistic Temporal Tensor Factorization* (P<sup>2</sup>T<sup>2</sup>F). First, we introduce a new data split strategy to divide the whole tensor data into several sub-tensors, and we meticulously reduce a large tensor operation to a set of independent operations toward sub-tensors allowing for concurrent executions. We also extend ADMM to handle these sub-tensors in the training.

##### A. Data Split Strategy

The costs of tensor decomposition are closely contingent upon tensor sizes, it is natural for us to split a large tensor data into several independent sub-tensors. In general, we divide the tensor  $\mathcal{X}$  into  $P$  sub-blocks (or sub-tensors) along the mode with the most dimensions (assuming the first mode without loss of the generality). In this case, each sub-block contains  $\frac{I}{P}$  horizontal slices (e.g.,  $\mathcal{X}^p \in \mathbb{R}^{\frac{I}{P} \times J \times K}$ ,  $p = 1, 2, \dots, P$ ).

Assuming there are  $P$  corresponding local factor matrices for each mode denoted as  $A^p$ ,  $B^p$  and  $C^p$ , respectively. Please note that  $B^p$  and  $C^p$  share the same size with  $B$  and  $C$ . To better utilize the global variable consensus optimization method, we have a global item latent matrix denoted as  $\bar{B}$  and a global time latent matrix denoted as  $\bar{C}$ . Since the local matrices  $B^p$  and  $C^p$  are only coupled with  $A^p$ , it is feasible to independently update  $A^p$ ,  $B^p$  and  $C^p$  for each process. This split strategy enables the CP decomposition problem to fit in the parallel ADMM framework.

Algorithm 1 demonstrates the details of proposed tensor data split strategy. Please note that it is also possible to divide a 3-order tensor into  $P$  sub-blocks simultaneously along two modes, yet the approach is subject to significant complex loss functions and constraints.

##### B. Parallel Probabilistic Temporal Tensor Factorization Model

In this split setting, the minimization problem (e.g., PTTF model) in (2) can be reformulated as the following constrained optimization problem:

$$\begin{aligned} \min_{\substack{A^p, B^p, C^p, \\ C_0^p, \bar{B}, \bar{C}}} \sum_{p=1}^P [f(A^p, B^p, C^p) + g(A^p, B^p, C^p, C_0^p)] \\ \text{s.t. } B^p - \bar{B} = \mathbf{0}, \\ C^p - \bar{C} = \mathbf{0}; \forall p \in \{1, 2, \dots, P\}. \end{aligned} \quad (3)$$

where

$$\begin{aligned} f(A^p, B^p, C^p) &= \frac{1}{2} \sum_{(i,j,k) \in \Omega^p} (x_{ijk}^p - \langle A_i^p, B_j^p, C_k^p \rangle)^2, \\ g(A^p, B^p, C^p, C_0^p) &= \frac{\lambda_A}{2} \|A^p\|_F^2 + \frac{\lambda_B}{2} \|B^p\|_F^2 \\ &\quad + \frac{\lambda_C}{2} \sum_{k=1}^K \|C_k^p - C_{k-1}^p\|_2^2 + \frac{\lambda_0}{2} \|C_0^p - \mu_C\|_2^2. \end{aligned}$$

Here,  $\Omega^p$  denotes the  $(i, j, k)$  indices of the values located in process  $p$ .  $\bar{B}$  and  $\bar{C}$  denote the global factor matrices. If we want to generalize it to a  $N$ -way tensor ( $N > 3$ ), we add additional constraints.

We transform the constrained optimization problem in (3) to an unconstrained problem with Augmented Lagrangian Method, and yield the following local objective function:

$$\begin{aligned} L^p(A^p, B^p, C^p, C_0^p, \Theta_B^p, \Theta_C^p, \bar{B}, \bar{C}) &= f(A^p, B^p, C^p) \\ &\quad + g(A^p, B^p, C^p, C_0^p) + l(B^p, C^p, \Theta_B^p, \Theta_C^p, \bar{B}, \bar{C}), \end{aligned} \quad (4)$$

where

$$\begin{aligned} l(B^p, C^p, \Theta_B^p, \Theta_C^p, \bar{B}, \bar{C}) &= \text{tr}([\Theta_B^p]^\top (B^p - \bar{B})) + (\rho_B/2) \|B^p - \bar{B}\|_F^2 \\ &\quad + \text{tr}([\Theta_C^p]^\top (C^p - \bar{C})) + (\rho_C/2) \|C^p - \bar{C}\|_F^2. \end{aligned}$$

Here,  $\Theta_B^p$  and  $\Theta_C^p$  denote the Lagrangian multipliers,  $\rho_B$  and  $\rho_C$  are the penalty parameters.

The global objective function is then as follows:

$$\begin{aligned} L(\mathcal{A}, \mathcal{B}, \mathcal{C}, \mathcal{C}_0, \Theta_B, \Theta_C, \bar{B}, \bar{C}) &= \sum_{p=1}^P L^p(A^p, B^p, C^p, C_0^p, \Theta_B^p, \Theta_C^p, \bar{B}, \bar{C}), \end{aligned} \quad (5)$$

where,  $\mathcal{A} = \{A^p\}_{p=1}^P$ ,  $\mathcal{B}, \mathcal{C}, \mathcal{C}_0, \Theta_B$  and  $\Theta_C$  are similarly defined. The ADMM method will solve this problem by repeating the following steps:

$$A_{t+1}^p, B_{t+1}^p, C_{t+1}^p, (C_0^p)_{t+1} \leftarrow \underset{A^p, B^p, C^p, C_0^p}{\text{argmin}} F_{\text{local}}, \quad (6)$$

$$\bar{B}_{t+1}, \bar{C}_{t+1} \leftarrow \underset{\bar{B}, \bar{C}}{\text{argmin}} F_{\text{global}}, \quad (7)$$

$$(\Theta_B^p)_{t+1} \leftarrow (\Theta_B^p)_t + \rho_B (B_{t+1}^p - \bar{B}_{t+1}), \quad (8)$$

$$(\Theta_C^p)_{t+1} \leftarrow (\Theta_C^p)_t + \rho_C (C_{t+1}^p - \bar{C}_{t+1}), \quad (9)$$

$$\forall p \in \{1, 2, \dots, P\},$$

where

$$F_{local} = L^p(\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p, C_0^p, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{\mathbf{B}}_t, \bar{\mathbf{C}}_t),$$

$$F_{global} = \sum_{p=1}^P l(\mathbf{B}_{t+1}^p, \mathbf{C}_{t+1}^p, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{\mathbf{B}}, \bar{\mathbf{C}}).$$

These update rules suggest that  $\mathbf{A}^p$ ,  $\mathbf{B}^p$ ,  $\mathbf{C}^p$ ,  $C_0^p$ ,  $\Theta_B^p$  and  $\Theta_C^p$  can be locally updated in an independent process. In this case, we dissect the whole tensor factorization problem into  $P$  independent sub-problems allowing for concurrent executions. Since we solve the problem under the ADMM framework, P<sup>2</sup>T<sup>2</sup>F can be viewed as a parallel extension of ADMM applied in PTTF problem.

### C. Stochastic Learning Algorithm for P<sup>2</sup>T<sup>2</sup>F

To learn the parameters in (3), we need to solve the above several steps from (6) to (9). If the optimal  $\mathbf{A}^p$ ,  $\mathbf{B}^p$  and  $\mathbf{C}^p$  have been obtained, it is easy to calculate  $C_0^p$ ,  $\bar{\mathbf{B}}$  and  $\bar{\mathbf{C}}$ . By setting the partial derivative of  $L^p$  w.r.t  $C_0^p$  to zero, we acquire the update rule of  $C_0^p$ :

$$C_0^p \leftarrow \frac{1}{\lambda_0 + \lambda_C} (\lambda_C C_1^p + \lambda_0 \mu_C). \quad (10)$$

Since  $\bar{\mathbf{B}}$  is a global variable, we need to take the partial derivative of  $L$  in (5) w.r.t  $\bar{\mathbf{B}}$ , and set the derivative to zero, then we can get  $\bar{\mathbf{B}}$ . If we set  $(\Theta_B^p)_0 = \mathbf{0}$ ,  $p = 1, \dots, P$ , we can prove that  $\sum_{p=1}^P (\Theta_B^p)_t = \mathbf{0}$ ,  $t = 1, 2, \dots$ . Then, the update rules for  $\bar{\mathbf{B}}$  and  $\bar{\mathbf{C}}$  ( $\bar{\mathbf{C}}$  is similar to  $\bar{\mathbf{B}}$ ) can be concisely written as:

$$\bar{\mathbf{B}} \leftarrow \frac{1}{P} \sum_{p=1}^P \mathbf{B}^p, \quad \bar{\mathbf{C}} \leftarrow \frac{1}{P} \sum_{p=1}^P \mathbf{C}^p. \quad (11)$$

$\Theta_B^p$  and  $\Theta_C^p$  can be directly updated by formulae (8) and (9). Therefore, how to efficiently compute the factor matrices becomes the key learning part. In the following content of this subsection, we will design a stochastic learning algorithm to solve it.

1) *Batch Learning*: The update step in (6) is actually a PTTF problem with  $\Theta_B^p$ ,  $\Theta_C^p$ ,  $\bar{\mathbf{B}}$  and  $\bar{\mathbf{C}}$  fixed. Since  $\mathbf{A}^p$ ,  $\mathbf{B}^p$  and  $\mathbf{C}^p$  are coupled together and the objective function of the PTTF problem is non-convex, it is not easy to get a satisfied solution. We employ a technique similar to that in [31] to resolve this issue by constructing a substitutive objective function, where the factor matrices  $\mathbf{A}^p$ ,  $\mathbf{B}^p$  and  $\mathbf{C}^p$  are decoupled and can be simultaneously calculated. The convexity of the constructed function, in each iteration, enables us to get the analytical solution of  $\mathbf{A}^p$ ,  $\mathbf{B}^p$  and  $\mathbf{C}^p$  by setting their gradients to zero.

The substitutive objective function is defined as follows:

$$H^p(\mathcal{M}^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{\mathbf{B}}_t, \bar{\mathbf{C}}_t, \tau_t | \mathcal{M}_t^p) \\ = h(\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p, \tau_t | \mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p) + g(\mathcal{M}^p, (C_0^p)_t) \\ + l(\mathbf{B}^p, \mathbf{C}^p, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{\mathbf{B}}_t, \bar{\mathbf{C}}_t), \quad (12)$$

where

$$h(\mathcal{M}^p, \tau_t | \mathcal{M}_t^p) = f(\mathcal{M}_t^p) + \text{tr}[\nabla_{\mathbf{A}^p}^\top f(\mathcal{M}_t^p)(\mathbf{A}^p - \mathbf{A}_t^p)] \\ + \text{tr}[\nabla_{\mathbf{B}^p}^\top f(\mathcal{M}_t^p)(\mathbf{B}^p - \mathbf{B}_t^p)] + \text{tr}[\nabla_{\mathbf{C}^p}^\top f(\mathcal{M}_t^p)(\mathbf{C}^p - \mathbf{C}_t^p)] \\ + \frac{1}{2\tau_t} (\|\mathbf{A}^p - \mathbf{A}_t^p\|_F^2 + \|\mathbf{B}^p - \mathbf{B}_t^p\|_F^2 + \|\mathbf{C}^p - \mathbf{C}_t^p\|_F^2). \quad (13)$$

Here,  $\mathcal{M}^p = \{\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p\}$ ,  $\mathcal{M}_t^p = \{\mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p\}$  are used for simple writing,  $\tau_t$  is a value that will be related to the learning rate and  $f(\mathcal{M}_t^p)$  is already defined in (3).

**Theorem 1.** Let  $D = \{\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p | \|\mathbf{A}_i^p - (\mathbf{A}_i^p)_t\|_2^2 \leq \delta^2, \|\mathbf{B}_j^p - (\mathbf{B}_j^p)_t\|_2^2 \leq \delta^2, \|\mathbf{C}_k^p - (\mathbf{C}_k^p)_t\|_2^2 \leq \delta^2\}$ ,  $\delta^2 > 0$ . Then,  $\forall \mathcal{M}^p \in D$ , a suitable  $\tau_t$  can always be found to make  $H^p(\cdot)$  satisfy the following two properties:

$$H^p(\mathcal{M}^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{\mathbf{B}}_t, \bar{\mathbf{C}}_t, \tau_t | \mathcal{M}_t^p) \\ \geq L^p(\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{\mathbf{B}}_t, \bar{\mathbf{C}}_t), \\ H^p(\mathcal{M}_t^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{\mathbf{B}}_t, \bar{\mathbf{C}}_t, \tau_t | \mathcal{M}_t^p) \\ = L^p(\mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{\mathbf{B}}_t, \bar{\mathbf{C}}_t).$$

The proof of Theorem 1 can be found in the supplemental material. From Theorem 1, we can find that  $H^p(\cdot)$  is an upper bound of  $L^p(\cdot)$ , and  $H^p(\cdot) = L^p(\cdot)$  at the point  $(\mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p)$ . Fortunately,  $H^p(\cdot)$  is convex in  $(\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p)$ , and  $\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p$  are decoupled in  $H^p(\cdot)$ . Hence, we can optimize the easily-solved constructed function  $H^p(\cdot)$  instead of  $L^p(\cdot)$  in (6) by setting the gradients to zero, the optimal results are computed as follows:

$$\mathbf{A}^p \leftarrow \frac{1}{1 + \lambda_A \tau_t} [\mathbf{A}_t^p - \tau_t * \nabla_{\mathbf{A}^p} f(\mathcal{M}_t^p)], \quad (14)$$

$$\mathbf{B}^p \leftarrow \frac{1}{1/\tau_t + \lambda_B + \rho_B} [\mathbf{B}_t^p / \tau_t + \rho_B \bar{\mathbf{B}}_t \\ - (\Theta_B^p)_t - \nabla_{\mathbf{B}^p} f(\mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p)], \quad (15)$$

$$\mathbf{C}^p \leftarrow \mathbf{Q}^{-1} [\mathbf{C}_t^p / \tau_t + \rho_C \bar{\mathbf{C}}_t + \lambda_C \mathbf{S}_C \\ - (\Theta_C^p)_t - \nabla_{\mathbf{C}^p} f(\mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p)], \quad (16)$$

where  $\mathbf{S}_C = \begin{bmatrix} (C_0^p)_t \\ \mathbf{0} \end{bmatrix}$  is a  $K$ -by- $R$  matrix and  $\mathbf{Q} = (1/\tau_t + \rho_C) \mathbf{I}_K + \lambda_C \mathbf{S}$  is a  $K$ -by- $K$  matrix. Here,  $\mathbf{S}$  denotes a coefficient matrix which shows relationships on the time dimension. Actually,  $\mathbf{S}$  is a tridiagonal matrix and can be described as follows:

$$\mathbf{S} = \begin{bmatrix} 2\mathbf{I}_{K-1} & \mathbf{0} \\ \mathbf{0} & 1 \end{bmatrix} + \begin{bmatrix} \mathbf{0} & -\mathbf{I}_{K-1} \\ 0 & \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{0} & 0 \\ -\mathbf{I}_{K-1} & \mathbf{0} \end{bmatrix}.$$

A batch learning algorithm for the problem in (3) can be got by combining (8)-(11) and (14)-(16).

**Theorem 2.** The batch learning algorithm enables P<sup>2</sup>T<sup>2</sup>F to converge to a local optimum.

*Proof.* From Theorem 1, we can get

$$\begin{aligned}
& L^p(\mathcal{M}_{t+1}^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{B}_t, \bar{C}_t) \\
& \leq H^p(\mathcal{M}_{t+1}^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{B}_t, \bar{C}_t, \tau_t | \mathcal{M}_t^p) \\
& \leq H^p(\mathcal{M}_t^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{B}_t, \bar{C}_t, \tau_t | \mathcal{M}_t^p) \\
& = L^p(A_t^p, B_t^p, C_t^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \bar{B}_t, \bar{C}_t).
\end{aligned}$$

That is to say, the global objective function  $L(\cdot)$  in (5) will not increase in each iteration. Furthermore,  $L(\cdot)$  is non-convex and has the lower bound  $-\frac{\sum_{p=1}^P \|\Theta_B^p\|_F^2}{2\rho_B} - \frac{\sum_{p=1}^P \|\Theta_C^p\|_F^2}{2\rho_C}$ . Hence, our batch learning algorithm will converge to a local optimum.  $\square$

2) *Stochastic Learning*: From the batch learning algorithm, we can find that the update rules for  $\mathcal{M}^p$  will need all ratings related to  $\mathcal{M}^p$ . If the number of ratings become very large, the batch learning algorithm will not be efficient. So, we propose a stochastic learning algorithm to further improve the efficiency, and the update rules for  $\mathcal{M}^p$  are as follows:

$$\begin{aligned}
A_i^p & \leftarrow \frac{(A_i^p)_t + \tau_t \epsilon_{ijk}((B_j^p)_t * (C_k^p)_t)}{1 + \lambda_A \tau_t}, \\
B_j^p & \leftarrow \frac{1}{1/\tau_t + \lambda_B + \rho_B} [(B_j^p)_t / \tau_t + \rho_B (\bar{B}_j)_t \\
& \quad - ((\Theta_B^p)_j)_t + \epsilon_{ijk}((A_i^p)_t * (C_k^p)_t)], \\
C_k^p & \leftarrow \frac{1}{1/\tau_t + 2\lambda_C + \rho_C} [(C_k^p)_t / \tau_t + \rho_C (\bar{C}_k)_t \\
& \quad + \lambda_C ((C_{k-1}^p)_t + (C_{k+1}^p)_t) \\
& \quad - ((\Theta_C^p)_k)_t + \epsilon_{ijk}((A_i^p)_t * (B_j^p)_t)], \quad (17)
\end{aligned}$$

where  $\epsilon_{ijk} = x_{ijk}^p - \langle (A_i^p)_t, (B_j^p)_t, (C_k^p)_t \rangle$ . Hence, the stochastic learning algorithm is a variant of the batch learning algorithm.

By combining the tensor data split strategy and the stochastic update rules stated above, we get a stochastic learning algorithm for our P<sup>2</sup>T<sup>2</sup>F model. The whole procedure of P<sup>2</sup>T<sup>2</sup>F is briefly listed in Algorithm 2, where  $\mathbf{A} = [(A^1)^\top (A^2)^\top \dots (A^P)^\top]^\top$ . Note that the convergence criterion is met when the difference between the train RMSEs of two successive iterations less than some threshold, e.g.,  $10^{-4}$ .

#### D. Complexity Analysis

P<sup>2</sup>T<sup>2</sup>F mainly needs two steps to update all variables once. The first step updates  $A^p$ ,  $B^p$  and  $C^p$ . For each value  $x_{ijk}$ , the time complexity of update  $A_i^p$ ,  $B_j^p$  and  $C_k^p$  is  $O(R)$ . Because the total number of observed entries in each process is about  $|\Omega|/P$ , the time complexity of step one is  $O(|\Omega|R/P)$ . The second step needs to update a matrix of size  $J \times R$  and a matrix of size  $K \times R$  in each process, so the time complexity is  $O(\max\{J, K\}R)$ . In total, the time complexity of P<sup>2</sup>T<sup>2</sup>F for each iteration can come down to  $O(|\Omega|R/P + \max\{J, K\}R)$ .

### V. EXPERIMENTAL RESULTS

Our experiments are designed to study the accuracy and efficiency of the P<sup>2</sup>T<sup>2</sup>F and baselines on the publicly available real-world datasets. All the experiments are run on a 12-core server with 2.60GHz Intel(R) Xeon(R) E5-2630 processor and 64GB of RAM.

#### Algorithm 2 Our P<sup>2</sup>T<sup>2</sup>F model

**Input:** Tensor  $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ , Rank  $R$ ,  $MaxIter$ ,  $P$ .  
**Output:**  $\mathbf{A} \in \mathbb{R}^{I \times R}$ ,  $\bar{\mathbf{B}} \in \mathbb{R}^{J \times R}$ ,  $\bar{\mathbf{C}} \in \mathbb{R}^{K \times R}$ .

- 1) use Algorithm 1 to get  $\mathcal{X}^p$ ;
- 2) initialize  $\lambda_A, \lambda_B, \lambda_C, \lambda_0, \rho_B, \rho_C, A^p, B^p, C^p, C_0^p, \mu_C$ ;
- 3) set  $\Theta_B^p$  (and  $\Theta_C^p$ ) =  $\mathbf{0}$ ,  $p = 1, 2, \dots, P$ ;
- 4) calculate  $\bar{B}, \bar{C}$  by (11);
- 5) **for**  $iter = 1, 2, \dots, MaxIter$  **do**
- 6)     **for**  $p = 1, 2, \dots, P$  **parallel do**
- 7)         update  $C_0^p$  by (10);
- 8)         **for** each  $x_{ijk}$  in process  $p$  **do**
- 9)             update  $A_i^p, B_j^p, C_k^p$  by (17);
- 10)         update  $\bar{B}, \bar{C}$  by (11);
- 11)         **for**  $p = 1, 2, \dots, P$  **parallel do**
- 12)             update  $\Theta_B^p, \Theta_C^p$  by (8) and (9);
- 13)         **if** convergence criterion is met **then**
- 14)             break;
- 15)         update  $\tau_t$ ;
- 16) **return**  $\mathbf{A}, \bar{\mathbf{B}}, \bar{\mathbf{C}}$ .

TABLE I  
SUMMARY OF REAL-WORLD DATASETS

	S1	S2	S3
$I_1$	14,012	28,060	56,361
$I_2$	19,527	24,981	28,444
$I_3$	242	242	242
#Train	1,851,291	3,739,047	7,566,903
#Test	205,699	415,449	840,766

a) *Datasets and Parameter Settings*: The real-world tensor data used in our experiments are public collaborative filtering datasets: Movielens ml-latest<sup>1</sup>, which is movie rating data from MovieLens, an online movie recommender service. In order to study P<sup>2</sup>T<sup>2</sup>F's parallel performance, we process it into three 3-order tensors, where each mode correspond to users, movies and calendar month, respectively, with the restriction of minimal 20 ratings per user. The rates range from 0.5 to 5, and the details are summarized in Table I. We set  $\lambda_A = \lambda_B = \lambda_C = 0.01$  and  $\rho_B = \rho_C = 0.5$ . Since it is difficult to compute the exact value for  $\tau_t$ , we approximately update it as  $\tau_{t+1} = \tau_t * \beta$  ( $0 < \beta < 1$ ) for the  $t$ -th iteration. We also set a threshold  $\alpha$ . When  $\tau_t \leq \alpha$ , we stop decreasing  $\tau_t$ . We set  $\tau_0 = 0.0005$ ,  $\beta = 0.9$  and  $\alpha = 0.0001$ . The hyperparameters are all determined by cross validation. We choose  $R = 20$  here.

b) *RMSE and Efficiency*: We use the root mean squared error (RMSE) to evaluate our P<sup>2</sup>T<sup>2</sup>F model and baselines. We first examine the significance of the improvement of P<sup>2</sup>T<sup>2</sup>F over the CP, PTTF and PPMF model on these datasets in one core by repeating the prediction tasks 12 times using different random initializations. Figure 2 shows the resulting

<sup>1</sup><http://grouplens.org/datasets/movielens/>

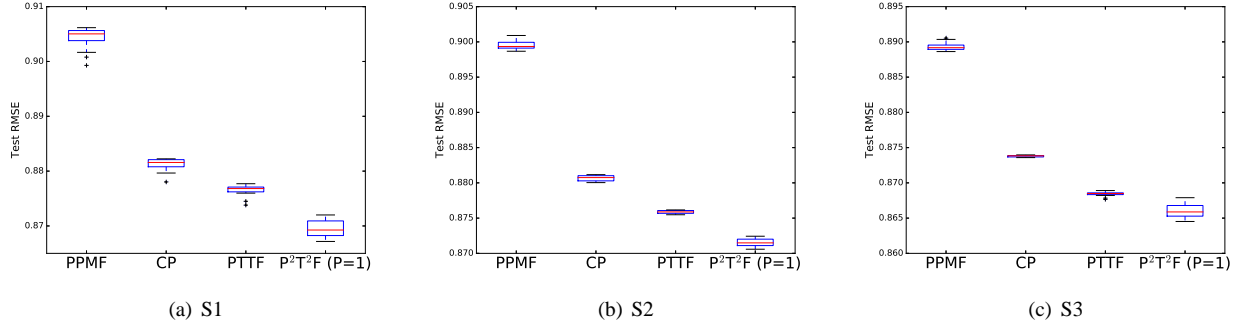


Fig. 2. Box plot of the RMSEs from PPMF, CP, PTTF and  $P^2T^2F$  on three datasets.  $P^2T^2F$  can outperform others in one core.

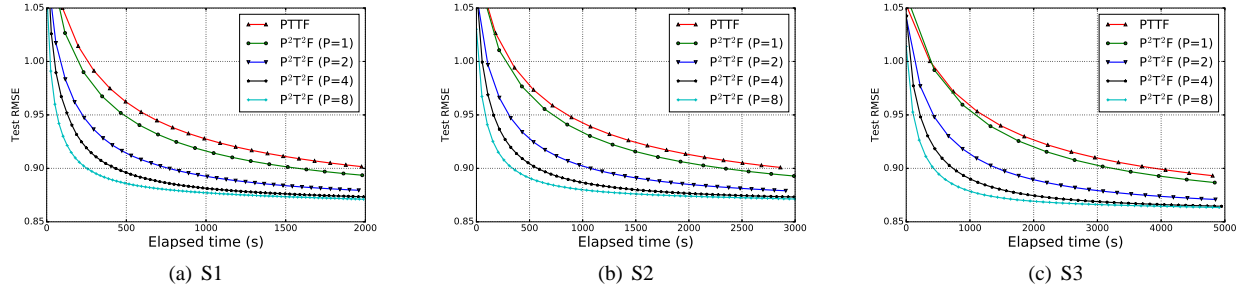


Fig. 3. The test RMSE curves for PTTF and  $P^2T^2F$  with different  $P$  on three datasets.

box plot of test RMSEs on S1-S3 at the moment when the convergence criterion of all these methods is satisfied. We can see that  $P^2T^2F$  model outperforms the baselines in all runs.  $P^2T^2F$  takes advantages of the stochastic learning algorithm, where it handles the latent factors using a surrogate objective function and makes them decoupled and easily computed. Therefore, it reaches a better RMSE than the conventional SGD-based PTTF method. In particular, the PTTF model outperforms CP model because of the temporal effects in the probabilistic decomposition. The PPMF model only considers the users and movies get the worst result. CP and PTTF model have a smaller degree of dispersion because they have fewer parameters than PPMF and  $P^2T^2F$  methods. Therefore, they are less sensitive random initializations.

*c) Scalability:* Another metric used to measure a parallel algorithm is the scalability. To study the scalability of  $P^2T^2F$ , we test our model on three datasets by varying the number of cores from 1 to 8. Figure 3 shows the test RMSE versus the running time for  $P^2T^2F$  model with different number of cores (or sub-tensors)  $P$ . The result demonstrates that the running time is approximately reduced to a half when the number of cores gets doubled. Note that the different curves of  $P^2T^2F$  eventually converge to the same solution. To see more clearly, we compute the speedup relative to the running time with 1 core ( $P = 1$ ) by varying the number of cores from 1 to 8. Here, we set  $RMSE = 0.90$  as a baseline. The results on S1-S3 are shown in Figure 4. We can intuitively see that  $P^2T^2F$  achieves nearly linear speedup and the speedup ratio increases with the number of ratings. The increased speedup from S1

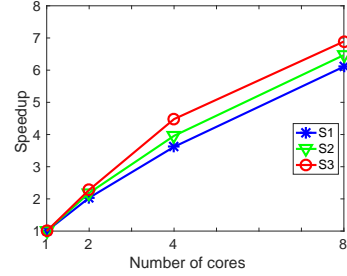


Fig. 4. The speedup of  $P^2T^2F$  w.r.t the number of cores on three datasets.

to S3 is probably caused by the increase of data density.

## VI. CONCLUSION

In this paper, we present  $P^2T^2F$  by deriving a stochastic ADMM algorithm to calculate the latent factors of probabilistic temporal tensors. We propose a new data split strategy to divide the large-scale problem into several independent sub-problems along the user dimension. Then we use the parallel ADMM framework to decompose these sub-tensors in parallel. Experiments on real world data sets demonstrate that our  $P^2T^2F$  model outperforms the traditional CP decomposition, PTTF and PPMF model in terms of efficiency and scalability.

## REFERENCES

- [1] L. De Lathauwer, J. Castaing, and J.-F. Cardoso, "Fourth-order cumulant-based blind identification of underdetermined mixtures," *IEEE Transactions on Signal Processing*, vol. 55, no. 6, pp. 2965–2973, 2007.

- [2] A. Shashua and A. Levin, "Linear image coding for regression and classification using the tensor-rank principle," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, vol. 1. IEEE, 2001, pp. 1–42.
- [3] G. Beylkin and M. J. Mohlenkamp, "Numerical operator calculus in higher dimensions," *Proceedings of the National Academy of Sciences*, vol. 99, no. 16, pp. 10 246–10 251, 2002.
- [4] Z. Xu, F. Yan, and Y. Qi, "Bayesian nonparametric models for multiway data analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 2, pp. 475–487, 2015.
- [5] S. Chen, M. R. Lyu, I. King, and Z. Xu, "Exact and stable recovery of pairwise interaction tensors," in *Advances in Neural Information Processing Systems*, 2013, pp. 1691–1699.
- [6] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.
- [7] L. R. Tucker, "Some mathematical notes on three-mode factor analysis," *Psychometrika*, vol. 31, no. 3, pp. 279–311, 1966.
- [8] L. Xiong, X. Chen, T.-K. Huang, J. G. Schneider, and J. G. Carbonell, "Temporal collaborative filtering with bayesian probabilistic tensor factorization," in *SDM*, vol. 10. SIAM, 2010, pp. 211–222.
- [9] A. Mnih and R. Salakhutdinov, "Probabilistic matrix factorization," in *Advances in neural information processing systems*, 2007, pp. 1257–1264.
- [10] R. Salakhutdinov and A. Mnih, "Bayesian probabilistic matrix factorization using markov chain monte carlo," in *Proceedings of the 25th international conference on Machine learning*. ACM, 2008, pp. 880–887.
- [11] I. Porteous, E. Bart, and M. Welling, "Multi-hdp: A non parametric bayesian model for tensor factorization," in *Aaai*, vol. 8, 2008, pp. 1487–1490.
- [12] M. N. Schmidt and S. Mohamed, "Probabilistic non-negative tensor factorization using markov chain monte carlo," in *Signal Processing Conference, 2009 17th European*. IEEE, 2009, pp. 1918–1922.
- [13] Y. Chi, S. Zhu, Y. Gong, and Y. Zhang, "Probabilistic polyadic factorization and its application to personalized recommendation," in *Proceedings of the 17th ACM conference on Information and knowledge management*. ACM, 2008, pp. 941–950.
- [14] Y. Koren, "Collaborative filtering with temporal dynamics," *Communications of the ACM*, vol. 53, no. 4, pp. 89–97, 2010.
- [15] Z.-Q. Yu, X.-J. Shi, L. Yan, and W.-J. Li, "Distributed stochastic admm for matrix factorization," in *Proceedings of the 23rd ACM Int. Conf. on Inform. and Knowl. Manag.* ACM, 2014, pp. 1259–1268.
- [16] B. Recht, C. Re, S. Wright, and F. Niu, "Hogwild: A lock-free approach to parallelizing stochastic gradient descent," in *Advances in Neural Information Processing Systems*, 2011, pp. 693–701.
- [17] R. Gemulla, E. Nijkamp, P. J. Haas, and Y. Sismanis, "Large-scale matrix factorization with distributed stochastic gradient descent," in *Proceedings of the 17th ACM conference on Knowledge discovery and data mining*. ACM, 2011, pp. 69–77.
- [18] L. Wang, W. Wu, Z. Xu, J. Xiao, and Y. Yang, "Blasx: A high performance level-3 blas library for heterogeneous multi-gpu computing," in *Proceedings of the 2016 International Conference on Supercomputing*. ACM, 2016, p. 20.
- [19] U. Kang, E. Papalexakis, A. Harpale, and C. Faloutsos, "Gigatensor: scaling tensor analysis up by 100 times-algorithms and discoveries," in *Proc. of the 18th ACM conf. on Knowl. disc. and dat. min.* ACM, 2012, pp. 316–324.
- [20] J. H. Choi and S. Vishwanathan, "Dfacto: Distributed factorization of tensors," in *Advances in Neural Information Processing Systems*, 2014, pp. 1296–1304.
- [21] A. P. Liavas and N. D. Sidiropoulos, "Parallel algorithms for constrained tensor factorization via alternating direction method of multipliers," *IEEE Transactions on Signal Processing*, vol. 63, no. 20, pp. 5450–5463, 2015.
- [22] A. H. Phan and A. Cichocki, "Parafac algorithms for large-scale problems," *Neurocomputing*, vol. 74, no. 11, pp. 1970–1984, 2011.
- [23] F. Shang, Y. Liu, and J. Cheng, "Generalized higher-order tensor decomposition via parallel admm," in *Proceedings of the 28th AAAI Conference on Artificial Intelligence*, 2014, pp. 1279–1285.
- [24] R. Tomioka, K. Hayashi, and H. Kashima, "Estimation of low-rank tensors via convex optimization," *ArXiv e-prints*, Oct. 2010.
- [25] Y. Wang, R. Chen, J. C. Denny, A. Kho, Y. Chen, B. A. Malin, and J. Sun, "Rubik: Knowledge guided tensor factorization and

completion for health data analytics," in *Proc. of the ACM Int. Conf. on Knowl. Disc. and Data Min.* ACM, 2015, pp. 1265–1274.

- [26] K. Shin and U. Kang, "Distributed methods for high-dimensional and large-scale tensor factorization," in *Data Mining (ICDM), 2014 IEEE International Conference on*. IEEE, 2014, pp. 989–994.
- [27] L. De Lathauwer, "A survey of tensor methods," in *Circuits and Systems, 2009. ISCAS 2009. IEEE International Symposium on*. IEEE, 2009, pp. 2773–2776.
- [28] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of COMPSTAT'2010*. Springer, 2010, pp. 177–186.
- [29] D. Goldfarb, S. Ma, and K. Scheinberg, "Fast alternating linearization methods for minimizing the sum of two convex functions," *Mathematical Programming*, vol. 141, no. 1-2, pp. 349–382, 2013.
- [30] J. Yang and Y. Zhang, "Alternating direction algorithms for  $\ell_1$ -problems in compressive sensing," *SIAM journal on scientific computing*, vol. 33, no. 1, pp. 250–278, 2011.
- [31] H. Ouyang, N. He, L. Tran, and A. G. Gray, "Stochastic alternating direction method of multipliers," *ICML (I)*, vol. 28, pp. 80–88, 2013.

## APPENDIX

We show the proof to Theorem 1 in the following.

### A. Proof of Theorem 1

*Proof.* The constructed function  $h(\cdot)$  in (13) can be written as

$$\begin{aligned} h(\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p, \tau_t | \mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p) \\ = \sum_{(i,j,k) \in \Omega^p} \hat{h}_{i,j,k}(\mathbf{A}_i^p, \mathbf{B}_j^p, \mathbf{C}_k^p, \tau_t | \mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p), \end{aligned}$$

where

$$\begin{aligned} \hat{h}_{i,j,k}(\mathbf{A}_i^p, \mathbf{B}_j^p, \mathbf{C}_k^p, \tau_t | \mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p) \\ = \hat{f}_{i,j,k}((\mathbf{A}_i^p)_t, (\mathbf{B}_j^p)_t, (\mathbf{C}_k^p)_t) \\ + \nabla_{\mathbf{A}_i^p} \hat{f}_{i,j,k}((\mathbf{A}_i^p)_t, (\mathbf{B}_j^p)_t, (\mathbf{C}_k^p)_t) (\mathbf{A}_i^p - (\mathbf{A}_i^p)_t)^T \\ + \nabla_{\mathbf{B}_j^p} \hat{f}_{i,j,k}((\mathbf{A}_i^p)_t, (\mathbf{B}_j^p)_t, (\mathbf{C}_k^p)_t) (\mathbf{B}_j^p - (\mathbf{B}_j^p)_t)^T \\ + \nabla_{\mathbf{C}_k^p} \hat{f}_{i,j,k}((\mathbf{A}_i^p)_t, (\mathbf{B}_j^p)_t, (\mathbf{C}_k^p)_t) (\mathbf{C}_k^p - (\mathbf{C}_k^p)_t)^T \\ + [1/(2m_i\tau_t)] \|\mathbf{A}_i^p - (\mathbf{A}_i^p)_t\|_2^2 \\ + [1/(2n_j\tau_t)] \|\mathbf{B}_j^p - (\mathbf{B}_j^p)_t\|_2^2 \\ + [1/(2z_k\tau_t)] \|\mathbf{C}_k^p - (\mathbf{C}_k^p)_t\|_2^2. \end{aligned}$$

Here,  $m_i$  denotes the number of ratings related to  $\mathbf{A}_i^p$  in  $\mathcal{X}^p$ ,  $n_j$  and  $z_k$  are similarly defined.

Then, we have

$$\begin{aligned} L^p(\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p, (\mathbf{C}_0^p)_t, (\boldsymbol{\Theta}_B^p)_t, (\boldsymbol{\Theta}_C^p)_t, \overline{\mathbf{B}}_t, \overline{\mathbf{C}}_t) \\ - H^p(\mathcal{M}^p, (\mathbf{C}_0^p)_t, (\boldsymbol{\Theta}_B^p)_t, (\boldsymbol{\Theta}_C^p)_t, \overline{\mathbf{B}}_t, \overline{\mathbf{C}}_t, \tau_t | \mathcal{M}_t^p) \\ = f(\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p) - h(\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p, \tau_t | \mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p) \\ = \sum_{(i,j,k) \in \Omega^p} [\hat{f}_{i,j,k}(\mathbf{A}_i^p, \mathbf{B}_j^p, \mathbf{C}_k^p) \\ - \hat{h}_{i,j,k}(\mathbf{A}_i^p, \mathbf{B}_j^p, \mathbf{C}_k^p, \tau_t | \mathbf{A}_t^p, \mathbf{B}_t^p, \mathbf{C}_t^p)]. \end{aligned}$$

For clarity, we denote  $A_i^p, B_j^p, C_k^p, (A_i^p)_t, (B_j^p)_t$  and  $(C_k^p)_t$  as  $\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{a}_t, \mathbf{b}_t$  and  $\mathbf{c}_t$ , respectively. Then we have

$$\begin{aligned}\hat{f}_{i,j,k}(\mathbf{a}, \mathbf{b}, \mathbf{c}) &= \frac{1}{2}(x_{ijk}^p - \langle \mathbf{a} - \mathbf{a}_t + \mathbf{a}_t, \\ &\quad \mathbf{b} - \mathbf{b}_t + \mathbf{b}_t, \mathbf{c} - \mathbf{c}_t + \mathbf{c}_t \rangle)^2 \\ &= \hat{f}_{i,j,k}(\mathbf{a}_t, \mathbf{b}_t, \mathbf{c}_t) \\ &\quad + \nabla_{\mathbf{a}} \hat{f}_{i,j,k}(\mathbf{a}_t, \mathbf{b}_t, \mathbf{c}_t)(\mathbf{a} - \mathbf{a}_t)^T \\ &\quad + \nabla_{\mathbf{b}} \hat{f}_{i,j,k}(\mathbf{a}_t, \mathbf{b}_t, \mathbf{c}_t)(\mathbf{b} - \mathbf{b}_t)^T \\ &\quad + \nabla_{\mathbf{c}} \hat{f}_{i,j,k}(\mathbf{a}_t, \mathbf{b}_t, \mathbf{c}_t)(\mathbf{c} - \mathbf{c}_t)^T \\ &\quad + o(\mathbf{a}, \mathbf{b}, \mathbf{c}),\end{aligned}$$

where  $o(\mathbf{a}, \mathbf{b}, \mathbf{c})$  contains all the second to sixth order terms.

We have the following properties by mainly using Cauchy inequality and the hypothesis  $\|\mathbf{a} - \mathbf{a}_t\|_2^2 \leq \delta^2, \|\mathbf{b} - \mathbf{b}_t\|_2^2 \leq \delta^2, \|\mathbf{c} - \mathbf{c}_t\|_2^2 \leq \delta^2$ :

$$\begin{aligned}-(x_{ijk}^p - \langle \mathbf{a}_t, \mathbf{b}_t, \mathbf{c}_t \rangle) &< \mathbf{a} - \mathbf{a}_t, \mathbf{b} - \mathbf{b}_t, \mathbf{c} - \mathbf{c}_t \rangle \\ &\leq |\epsilon_{ijk}|(\mathbf{a} - \mathbf{a}_t)((\mathbf{b} - \mathbf{b}_t) * \mathbf{c}_t)^T \\ &\leq \frac{1}{2}|\epsilon_{ijk}|(\|\mathbf{a} - \mathbf{a}_t\|_2^2 + \|\mathbf{b} - \mathbf{b}_t\|_2^2 \|\mathbf{c}_t\|_2^2); \\ -(x_{ijk}^p - \langle \mathbf{a}_t, \mathbf{b}_t, \mathbf{c}_t \rangle) &< \mathbf{a} - \mathbf{a}_t, \mathbf{b} - \mathbf{b}_t, \mathbf{c} - \mathbf{c}_t \rangle \\ &\leq |\epsilon_{ijk}|(\mathbf{a} - \mathbf{a}_t)((\mathbf{b} - \mathbf{b}_t) * (\mathbf{c} - \mathbf{c}_t))^T \\ &\leq \frac{1}{2}|\epsilon_{ijk}|(\|\mathbf{a} - \mathbf{a}_t\|_2^2 + \|\mathbf{b} - \mathbf{b}_t\|_2^2 \|\mathbf{c} - \mathbf{c}_t\|_2^2) \\ &\leq \frac{1}{2}|\epsilon_{ijk}|(\|\mathbf{a} - \mathbf{a}_t\|_2^2 + \|\mathbf{b} - \mathbf{b}_t\|_2^2 \delta^2);\end{aligned}$$

where  $\epsilon_{ijk} = x_{ijk}^p - \langle \mathbf{a}_t, \mathbf{b}_t, \mathbf{c}_t \rangle$ .

$$(a_1 + a_2 + \dots + a_n)^2 \leq n(a_1^2 + a_2^2 + \dots + a_n^2);$$

$$\begin{aligned}(\langle \mathbf{a} - \mathbf{a}_t, \mathbf{b} - \mathbf{b}_t, \mathbf{c} - \mathbf{c}_t \rangle)^2 &\leq \frac{1}{4}(\|\mathbf{a} - \mathbf{a}_t\|_2^2 + \|\mathbf{b} - \mathbf{b}_t\|_2^2 \delta^2)^2 \\ &\leq \frac{1}{4}(\delta^2 + \delta^4)(\|\mathbf{a} - \mathbf{a}_t\|_2^2 + \|\mathbf{b} - \mathbf{b}_t\|_2^2 \delta^2);\end{aligned}$$

$$\begin{aligned}(\langle \mathbf{a} - \mathbf{a}_t, \mathbf{b} - \mathbf{b}_t, \mathbf{c}_t \rangle)^2 &\leq \frac{1}{4}(\|\mathbf{a} - \mathbf{a}_t\|_2^2 + \|\mathbf{b} - \mathbf{b}_t\|_2^2 \|\mathbf{c}_t\|_2^2)^2 \\ &\leq \frac{1}{4}(\delta^2 + \delta^2 \|\mathbf{c}_t\|_2^2)(\|\mathbf{a} - \mathbf{a}_t\|_2^2 + \|\mathbf{b} - \mathbf{b}_t\|_2^2 \|\mathbf{c}_t\|_2^2); \\ (\langle \mathbf{a} - \mathbf{a}_t, \mathbf{b}_t, \mathbf{c}_t \rangle)^2 &\leq (\|\mathbf{a} - \mathbf{a}_t\|_2 \|\mathbf{b}_t * \mathbf{c}_t\|_2)^2 \\ &= \|\mathbf{b}_t * \mathbf{c}_t\|_2^2 \|\mathbf{a} - \mathbf{a}_t\|_2^2.\end{aligned}$$

Using the above six properties, we can prove that

$$\begin{aligned}o(\mathbf{a}, \mathbf{b}, \mathbf{c}) &\leq \pi_A \|\mathbf{a} - \mathbf{a}_t\|_2^2 + \pi_B \|\mathbf{b} - \mathbf{b}_t\|_2^2 \\ &\quad + \pi_C \|\mathbf{c} - \mathbf{c}_t\|_2^2,\end{aligned}$$

where  $\pi_A, \pi_B, \pi_C$  are constants which depend on  $\mathbf{a}_t, \mathbf{b}_t, \mathbf{c}_t$  and  $\delta^2$ .

If we let

$$\frac{1}{\tau_t} \geq \max\{2m_i \pi_A, 2n_j \pi_B, 2z_k \pi_C\},$$

then we can prove that

$$\begin{aligned}H^p(\mathcal{M}^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \overline{\mathbf{B}}_t, \overline{\mathbf{C}}_t, \tau_t | \mathcal{M}_t^p) \\ \geq L^p(\mathbf{A}^p, \mathbf{B}^p, \mathbf{C}^p, (C_0^p)_t, (\Theta_B^p)_t, (\Theta_C^p)_t, \overline{\mathbf{B}}_t, \overline{\mathbf{C}}_t).\end{aligned}$$

The second property in Theorem 1 can be easily proved.  $\square$